# Benchmarking, Monitoring, Observability, and Experimenting for Big Data and Machine Learning Systems

*Hong-Linh Truong*
*Department of Computer Science*
*linh.truong@aalto.fi, https://rdsea.github.io*

CS-E4660 Advanced Topics in Software Systems, Fall 2023
13/09/2023

# Learning objectives

- **Able to analyze the role of measurement, monitoring and observability in real-world cases for R3E**

- **Understand and develop methods with key steps and important tools for benchmarking, monitoring, observability and experimenting**

- **Able to apply these methods for big data/ML systems**

# The role of measurement, monitoring and observability

**Aalto University**
**School of Science**

# Development vs Runtime activities

**Design, test and benchmark R3E**

- **R3E for individual components**
- **model/capture complex dependencies**
- **design logs, metrics and traces for capturing states and complex dependencies**
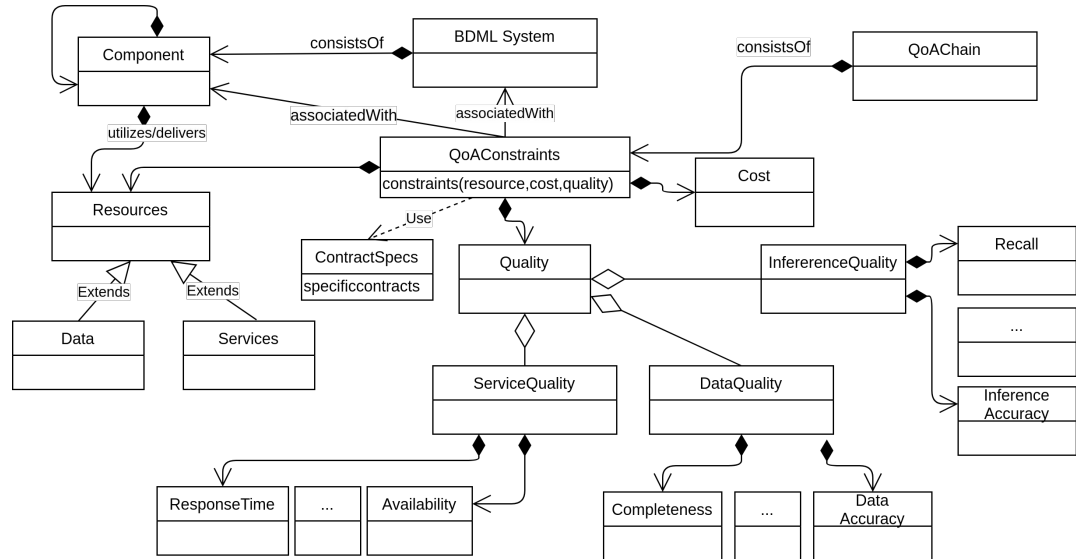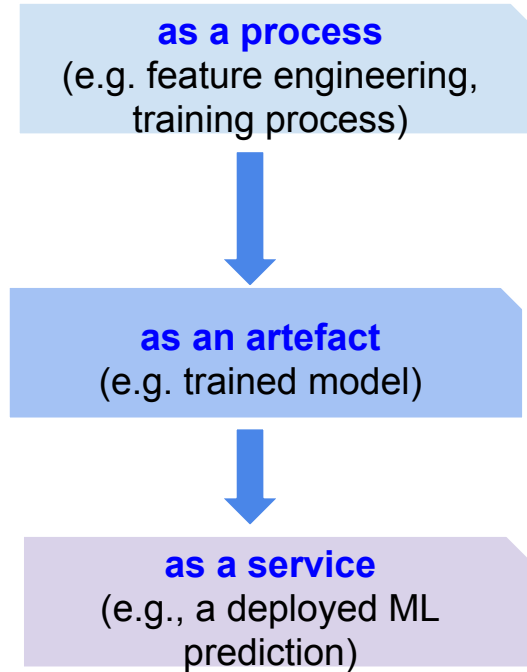
**Monitoring/observability and runtime adaptation**

- **runtime monitoring and observability**
- **states, performance and failure analytics**
- **runtime controls (constraints, rules, actions)**

# Measurement, monitoring, and observability for R3E

- **Instrumentation and sampling**
  - instrumentation: insert <span style="color:red">probes into systems</span> to measure system behaviors directly or produce logs
  - sampling: use components to sample system behaviors
- **Monitoring**
  - perform sampling or instrumentation to collect and share metrics, logs, traces; visualize what has been happened
- **Observability**
  - evaluate and interpret measurements for specific contexts
  - understand and explain the systems states, dependencies, etc.

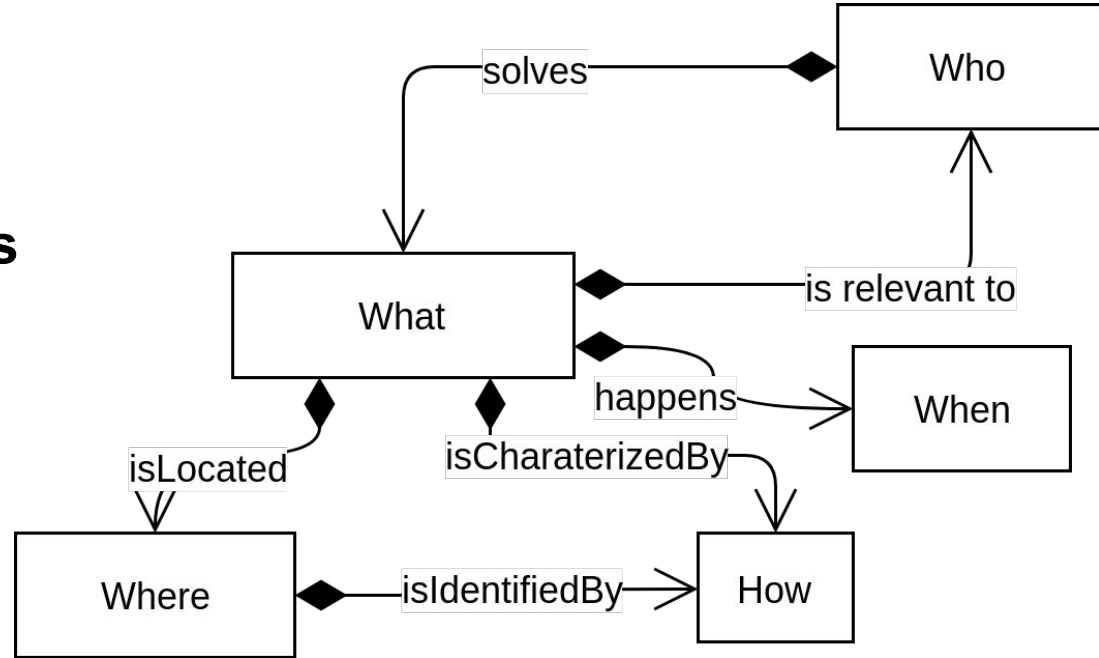# Recall: strongly interdependencies

**Any problem would lead to a huge waste (engineering effort, operation cost, societal impact due to wrong inference/prediction)**

**as a process**
(e.g. feature engineering, training process)

**as an artefact**
(e.g. trained model)

**as a service**
(e.g., a deployed ML prediction)

# Methods

**Aalto University**
**School of Science**

# What/Which, Where, When, Who and How

**Understand W4H aspects for analytics of big data/ML systems**

# Key steps – What/Which

- **Understand and identify indicators/metrics characterizing your systems**
- **Common metrics vs specific (big data/ML) ones**
  - different relevance/importance based on specific contexts
- **Most critical problems are due to complex dependencies that are not common**
  - root cause analysis will be tricky
- **For which purposes?**
  - SRE, benchmarking, Test-Driven Development (TDD)

# Key steps – Where and When

- **Where: as a " space" dimension**

  - tightly coupled or isolated/loosely coupled

  - different places

    - software/system layers, components and systems boundaries

    - dependencies among components

    - development/configuration pipelines

- **When: as a „time" dimension**

  - design, test/training, or runtime (DevOps)

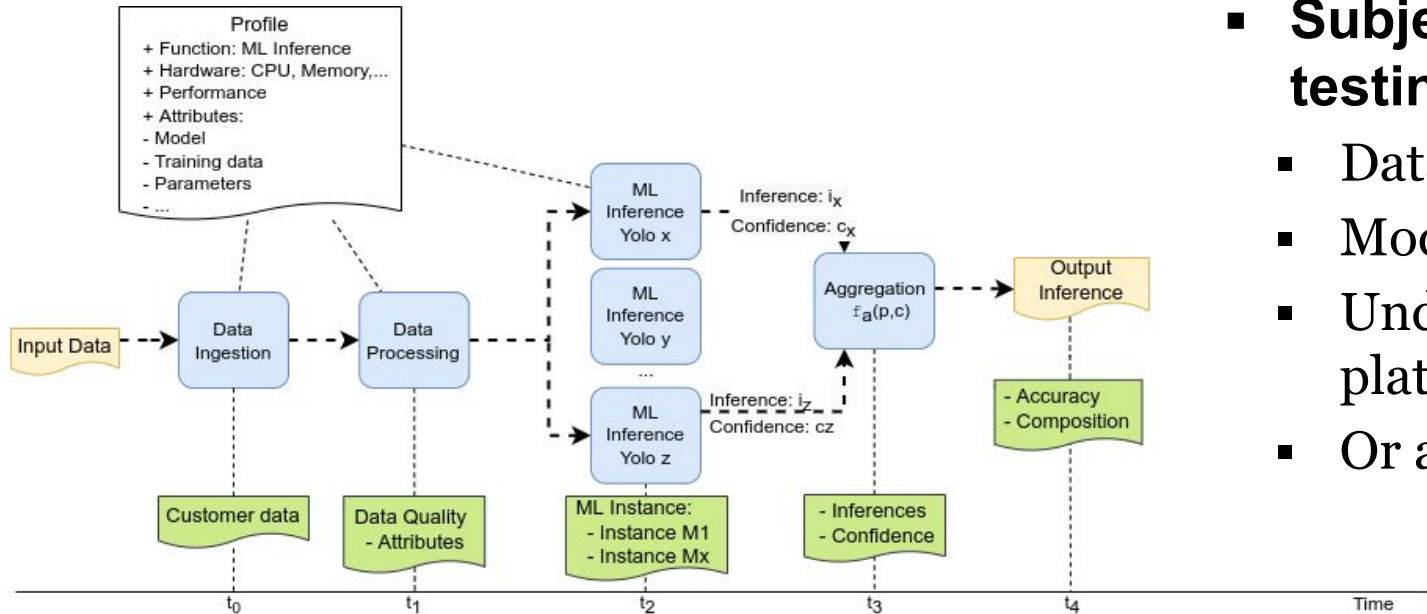  - further divided into sub states

# Key steps - How

- **Characterize dependencies among components**
  - understand the system as a whole
  - include also development processes, data, software artefacts and execution environments
- **Select tools for capturing metrics**
- **Understand what kind of changes/designs we must do**
- **Do monitoring and analysis**
- **Integrate many types of data for monitoring and observability**

# Apply W4H for benchmarking, monitoring, validation and experimenting

- **Determines clearly system boundaries**
    - the system under study, the system used to judge, and the environment
    - "domain-driven/oriented" and bounded context principles
- **Understands dependencies**
    - among components in distributed big data/ML systems in distributed computing platforms
    - single layer as well as cross-layered dependencies
- **Determines types of metrics and failures and break down problems along the dependency path (how)**
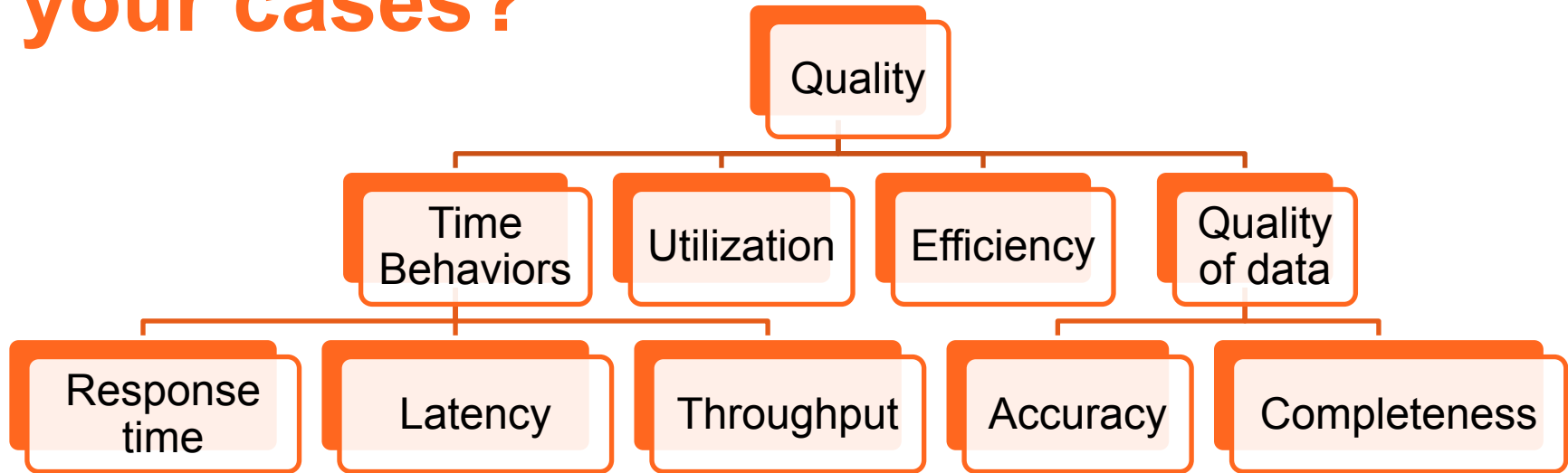
# Boundaries and dependencies

**Example of a ML service for object recognition (used in our hands-on)**



- **Subjects for testing/debugging**
  - Data?
  - Model?
  - Underlying service platform?
  - Or all of them?

Aalto University
School of Science

# What are the most critical metrics for your cases?

Quality

- Time Behaviors
- Utilization
- Efficiency
- Quality of data

- Response time
- Latency
- Throughput
- Accuracy
- Completeness

Industry view: https://guidingmetrics.com/content/cloud-services-industrys-10-most-critical-metrics/
NIST:  https://www.nist.gov/sites/default/files/documents/itl/cloud/RATAX-CloudServiceMetricsDescription-DRAFT-20141111.pdf

**Contradiction/Tradeoffs between Efficiency versus Resiliency Metrics for an ML model =! Metrics for ML system**

# Common performance metrics

- **Timing behaviors**
  - Communication
    - *Latency/Transfer time*
    - *Data transfer rate, bandwidth*
  - Processing
    - *Response time (service latency/time)*
    - *Throughput*
- **Utilization**
  - Network utilization
  - CPU utilization
  - Service utilization
- **Efficiency/Scalability**
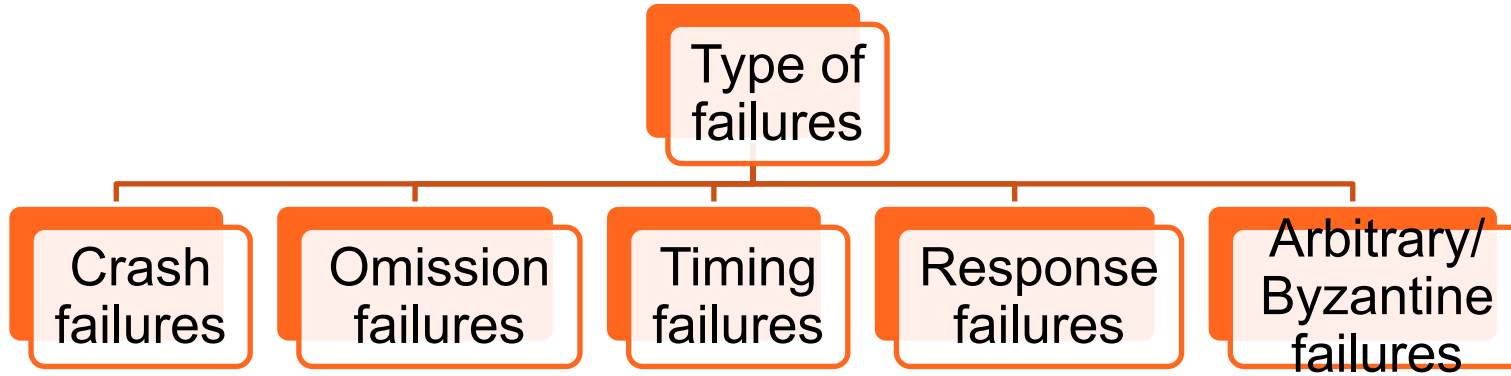  - Concurrent executions

## are they enough for big data/ML systems?

**Client**    **ML Service**

request

latency

Sending time

request

Processing time (service latency)

End-to-end response time

Receiving time

# Types of Failure

**Common**

```
              ┌──────────────┐
              │   Type of    │
              │   failures   │
              └──────┬───────┘
   ┌──────┬──────────┼──────────┬──────────┐
┌──────┐┌──────┐┌──────┐┌──────────┐┌──────────┐
│Crash ││Omission││Timing││Response  ││Arbitrary/│
│failures││failures││failures││failures ││Byzantine │
└──────┘└──────┘└──────┘└──────────┘│failures  │
                                      └──────────┘
```

**But unforeseen failures cannot be determined in advance ☐ design for handling failure**

**Check: https://arxiv.org/pdf/1910.11015.pdf for a "Taxonomy of Real Faults in Deep Learning Systems"**

# Metrics for Data

- **Completeness**
- **Timeliness**
- **Currency**
- **Validity**
- **Format**
- **Accuracy**
- **Data Drift**

**Understand the impact** →

**Forecasting**



**Drift impact**



**(examples with real mobile data)**

Often evaluation methods are different for different types of data

Aalto University
School of Science

# Metrics for ML models

- **Confusion matrix**

- **Accuracy**

- **Loss**

- **True positive rate**

- **False positive rate**

- **F1 Score/F-measure**

- **Etc.**

**(see https://towardsdatascience.com/metrics-to-evaluate-your-machine-learning-algorithm-f10ba6e38234)**

**How would we define "reliable function" of the model? E.g., when should we "retrain" the model?**

But how can we relate them to ML service/system?

**Aalto University
School of Science**

# Benchmarking and Observability

# Benchmarking

- **Benchmark: for comparing big data/ML systems w.r.t. selected (standard/common) workloads**

- **Where to be benchmarked**

  - benchmark individual subsystems: message brokers and data ingestion, databases and ingestion/query, data processing, ML models, serving platform

- **What to be benchmarked**

  - data ingestion throughput, processing throughput and time, component CPU and memory

  - training and inferencing time and accuracy

# Benchmarking

## What should we do for a big data system?



Check:
https://www.sciencedirect.com/science/article/pii/S0140366419312344
https://www.benchcouncil.org/BigDataBench/

# Benchmarking

**If you have an end-to-end ML system, does it make sense to benchmark the whole system?**

**Aalto University
School of Science**

# Benchmarking - ML

**Examples:**



Source: https://mlcommons.org/en/inference-edge-31/

**Also check: https://www.benchcouncil.org/aibench/index.html**
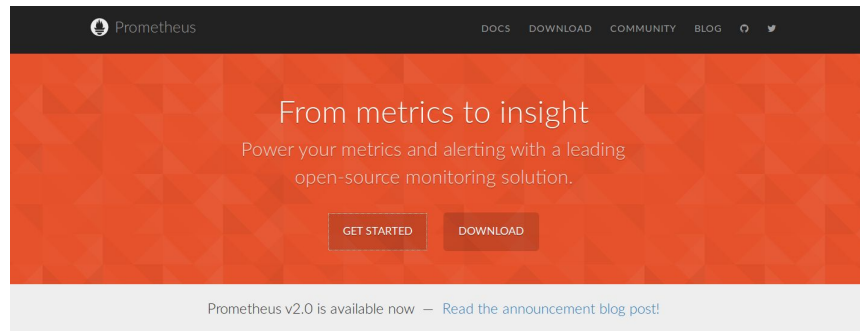
**Aalto University
School of Science**

# Service/Infrastructure monitoring tools

**There are many powerful tools!**

**But only low-level, well-identified monitoring data (infrastructures): pre-defined metrics exposed through interfaces with push/pull mechanism**

From: https://prometheus.io/

Aalto University
School of Science

# Instrumentation for observability

**Code instrumentation: for many metrics and logs that cannot be obtained from the outside of the component**

**the developer can instrument the code to capture metrics/generate logs/traces**



From: https://www.fluentd.org/



Filebeat

Lightweight shipper for logs

https://www.elastic.co/beats/filebeat



https://opentelemetry.io/

Aalto University
School of Science

# Monitoring data metrics on-the-fly

```
Data sources          Messaging systems        Stream processing
(sensors, files,  -->  (e.g., Kafka, AMQP, -->  systems
database, queues, log  MQTT, Pulsar)            (e.g. Flink, Kafka,
services)                                       Google Dataflow)
```

Operation/Management/Business Services

Dynamic ML Prediction systems (e.g.Spark, TensorFlow)

**Data drift detection? Quality of data detection? Or performance prediction?**

# A couple of tools for data quality

- **Generic tools/framework for checking data at rest**
  - Great expectation: https://github.com/great-expectations/great_expectations
  - YData (https://github.com/ydataai/ydata-quality)
  - Alibi-Detect (https://github.com/SeldonIO/alibi-detect)
  - Why-log (https://docs.whylabs.ai/docs/whylogs-overview/)
- **Integrated with processes in specific systems**
  - https://aws.amazon.com/blogs/industries/how-to-architect-data-quality-on-the-aws-cloud/
- **Working with specific data processing frameworks**
  - https://github.com/awslabs/python-deequ

# Visualization

**Metrics and Visualization**

- **Easy to visualize many types of metrics**
  - Human-in-the-loop
- **But only you can specify, define and map them to your structured applications**
- **Not for complex process automation!**
  - further integration and intelligence analytics (ML?)



https://www.elastic.co/products/kibana



https://grafana.com/

Aalto University
School of Science

# Observability

- **To monitor and understand the system as whole, end-to-end**
  - every component must be monitored
  - dependencies/interactions must be captured
  - diverse metrics, logs, tracing, etc. are needed to be integrated
- **Understand the states and behaviors of the whole systems**
- **Complex problems in big data/ML systems as these systems**
  - large-scale number of microservices in large-scale virtualized infrastructures
  - multi-dimensional states (code, models and data)

# Understand the structure of big data/ML application

**Dependency Structure**

- **Composable method**
  - divide a complex structure into basic common structures
  - each basic structure has different ways to analyze specific failures/metrics
- **Interpretation based on context/view**
  - client view or service provider view?
  - conformity versus specific requirement assessment



Client: Service is failed

Provider: OK

Failure

Slow

Aalto University
School of Science

# Support an end-to-end view or not

- **End-to-end reflects the entire system**
  - e.g., data reliability: from sensors to the final analytics/inference results
  - what if the developer/provider cannot support end-to-end?
- **The user expects end-to-end R3E**
  - e.g., specified in the expected accuracy
- **Providers/operators want to guarantee end-to-end quality**
  - need to monitor different parts, each has subsystems/components
  - coordination-aware assurance, e.g., using elasticity

# Big data/ML for Observability vs Observability for Big data/ML systems

- **Big data of metrics, logs and traces**
  - Large number of entities to be observed
  - High number of measurement dimensions
- **ML for observability**
  - Classification, prediction and detection of traffics/interactions anomaly behaviors, hidden relationships, etc.
  - Root-cause analysis
  - ML serving is in the edge and cloud

# Example: ML contract observability with QoA4ML



**https://github.com/rdsea/QoA4ML**

# Experiment management

# how do we manage important information for ML services?

# Problems

- **We need to run many experiments**
  - testability/observability purposes: figure out suitable configurations
  - how does this help to understand and support R3E?

- **Experiment management**
  - known domain and well-known books (e.g., "Design and Analysis of Experiments" by Douglas C. Montgomery)
  - principles: capturing various configurations
  - how does it work in big data and ML?

- **What do we need?**
  - tools/frameworks for tracking experiments

# Notions

- **A single run/trial**
  - inputs, results, required software artefacts
  - computing resources, logs/metrics
- **Experiment**
  - a collection of runs/trials/executions gathered in a specific context
- **Steps**
  - parameterization: generate different parameters
  - deployment: prepare suitable environments
  - execution: run and collect metrics
  - analysis and sharing: analyze experiment data

# Experiment tracking



Pipeline1:

doDataCollection:
    …
    **start(concern1)**
    getData(…)
    **stop(concern1)**
    …
    …

Pipeline2:

doFeatureExtraction:
    …
    **start(concern2)**
    extractFeature()
    **stop(concern2)**
    …

Pipeline3:

doClassification:
    …

    #do something;

Probe: instrumented code for logs, metrics, etc

Monitoring/ Experiment Platform

Experiment Data (Database, Server, Files)

Points of instrumentation

**But remember it is very large system! Different techniques/tools may be needed**

# Examples

- **Tensorflow Board (https://www.tensorflow.org/tensorboard)**
- **Experiment in Azure ML SDK**
  - https://docs.microsoft.com/en-us/python/api/overview/azure/ml/?view=azure-ml-py#experiment
- **MLFlows  https://mlflow.org/**
- **Kubeflows**
  - https://www.kubeflow.org/docs/pipelines/overview/concepts/
- **DVC:  https://dvc.org/**
- **Verta: https://www.verta.ai/**
- **Comet: https://www.comet.com/**

# Examples: MLFlow APIs

- **Experiment**

```
mflow.start_run()/end_run()

mflow.autolog()
```

- **Logs/metrics collection**

```
mflow.set_tag()

mflow.log_*()
```

- **Tracking data management**
    - Local files, Databases, HTTP server, Databrick logs

**(follow our hands-on tutorial)**

# Experiment management: more than just ML models

- **Remember there are many components in a system**

- **Experiment data about other components is also crucial**
    - have a full visibility and understanding of the system
    - support explainability and end-to-end optimization

- **ML model experiment must be combined with other types of experimental data**
    - experiment management for end-to-end systems

# Study log 2

**Describe one big data/ML pipeline that you are familiar with and explain your thoughts on how would you support the aspects of "benchmarking", "monitoring", "observability",  or "experimenting" for testing/implementing R3E aspects**

- Is enough to focus on 1 pipeline and 1 aspect
  - *No "familiar pipeline" → look at our hands-on tutorials*
- Be concrete, e.g., with metrics and possible tools
- Analyze if things can be done easily or where are the challenges that might be interesting for further investigation
- Optionally link to issues raised/addressed in a reading paper

**Aalto University**
**School of Science**

# Thanks!

**Hong-Linh Truong**
**Department of Computer Science**

**rdsea.github.io**